

UNITED STATES PATENT APPLICATION

INVENTORS:

JOHN M. QUERNEMOEN
MARK G. HAZZARD

APPLICATION:

BUILT IN HEADROOM FOR AN NT SYSTEM SIZER

ATTORNEY DOCKET NO.

RA-5247

Charles A. Johnson
Attorney for Applicant
Reg. No. 20,852
Telephone No. (651) 635-7702
Unisys Corporation
M.S. 4773
PO Box 64942
St. Paul, MN 55164-0942

CERTIFICATE UNDER 37 CFR 1.10: The undersigned hereby certifies that this transmittal letter and the paper of papers, as described hereinabove, are being deposited in the United States Postal Service, "Express Mail Post Office to Addressee" having an Express Mail mailing label number of **EK084997203US**, in an envelope addressed to: Box PATENT APPLICATION, Assistant Commissioner for Patents, Washington, D.C. 20231 on this 29th day of February, 2000.

Charles A. Johnson
(Charles A. Johnson)

February 29, 2000
(Date)

Related Applications

IN.S. B17
The present application is related to U.S. Patent Application Serial No. _____, filed _____, entitled DATABASE SIZER FOR NT SIZER
5 SYSTEM; U.S. Patent Application Serial No. _____, filed _____, entitled
SIZING SERVERS FOR DATABASE MANAGEMENT SYSTEMS VIA USER
DEFINED WORKLOADS; U.S. Patent Application Serial No. _____, filed
_____, entitled COMBINATION OF MASS STORAGE SIZER, COMPARATOR,
10 OLTP USER DEFINED WORKLOAD SIZER, AND DESIGN TRADE-OFF TOOL
IN ONE PACKAGE; U.S. Patent Application Serial No. _____, filed _____,
entitled ALGORITHMS TO CALCULATE MASS STORAGE REQUIREMENTS
FOR NT SIZER; and U.S. Patent Application Serial No. _____, filed
_____, entitled METHOD OF COMPARISON FOR COMPUTER SYSTEMS AND
15 APPARATUS THEREFOR, all of which are assigned to the assignee of the present
invention and incorporated herein by reference.

Field of the Invention

The present invention relates generally to computers and software. More
specifically, the present invention is related to database management system (DBMS)
hardware and software. The present invention includes software tools for determining
20 the hardware components required to handle given workloads while adjusting the
components to remain within hardware utilization limits.

Background of the Invention

Relational databases came into common use in computers over twenty years
ago. Despite improvements in database software and new methodologies, relational
25 databases remain the mainstay of database management systems (DBMS). Hardware

vendors originally supported proprietary database management systems which ran primarily on machines manufactured by the hardware vendor. Software developers later developed database management systems that were more open and ran on computers made by several vendors. The database management systems were also
5 ported to run under various operating systems. This gave the advantage of spreading the cost of development over more sites and also uncoupled the dependence between hardware vendors and software vendors. Third party support and training also became more common.

Database management systems also became separated into client-side software
10 and server-side software. This meant that the server-side software was decoupled from software having to do with the display, use, and formatting of the data received from the database. In particular, server-side software often handled mostly queries of existing data along with updates of existing data and insertion of new data.

Modern electronic commerce such as commerce over the Internet or business-
15 to-business electronic commerce has placed increased demands on many servers. This has also made frequent upgrades necessary. Company mergers and acquisitions frequently make it necessary to incorporate large amounts of data from unexpected sources. Customer expectations also make it necessary to upgrade hardware to keep up with the faster response times users expect even though system loads may be
20 increasing as well.

When installing, upgrading or replacing database servers it is desirable to establish the transaction handling requirements that must be met, and compare these requirements with the DBMS transaction handling processing capabilities of available systems.

25 In some situations, a required or desired transaction handling capability is

based mainly on the capability of a known system. It may be the case that a given brand name DBMS server is believed to satisfy a current or future requirement, and the tpmC capability of that DBMS server is available from the vendor and/or from the tpmC organization database. A more accurate performance value could, in theory, be derived from a series of more specific description of the requirements. It may be the case that the user has a more specific description of the requirements for a system, such as detailed transaction specific information.

When evaluating such systems, it would be desirable to allow for some margin or headroom in capacity. If a system was specified to run too close to capacity, then it might become backlogged and deny service when bursts of activity exceeded the average workload. It would also be desirable to factor in hardware utilization limits for specific hardware components. An upper utilization limit could reduce the likelihood of under-capacity and the resulting bottleneck. A lower utilization limit could reduce the likelihood of over-capacity and the resulting excessive spending.

What would be desirable therefore is a system for factoring in hardware utilization limits in determining hardware requirements. This may include allowing for changes in the hardware utilization limits, with corresponding adjustments to the hardware requirements.

Summary of the Invention

The present invention includes methods and systems for determining the hardware requirements for a database management system. In particular, the present invention includes methods and systems for determining the hardware required for an on-line transaction processing system. The hardware requirements are calculated so as to remain within user supplied hardware utilization limits.

In one illustrative embodiment, the server type, maximum processor

utilization, maximum network interface card utilization, and expected workload transactions per minute requirements are input to the system. The method calculates the required number of processors and the speed of network interface cards. In some methods, the database requirements are also input to the system and the required amount of memory and data storage are also calculated as part of the same program.

In one program according to the present invention, the program initializes hardware resource utilization limits to default hardware utilization limits. The program obtains throughput workload requirements from a human user and also obtains database requirements from the user. The program then calculates the hardware resources required to satisfy the database management system database and transaction handling requirements. The calculation determines the hardware requirements so as to remain within the user supplied utilization limits. The program then displays or otherwise outputs the hardware requirements to the user. The program may then accept user changes either to the workload requirements or to the hardware utilization limits. The hardware resources required can then be recalculated, again within the hardware resource utilization limits. In one method, the hardware resource utilization limits include processor utilization and network interface card utilization limits. The hardware resource utilization limits preferably include both upper and lower limits. The upper limits can serve to protect against bottlenecks occurring during bursts of activity. The lower limits can serve to avoid purchasing more hardware than needed. Using the present invention, the suggested number of discrete hardware components such as processors, memory, disk drives, and NICs can be automatically changed in response to changing requirements and utilization limits.

Brief Description of the Drawings

Figure 1 is a highly diagrammatic schematic of a computer system including a

database server;

Figure 2 is a highly diagrammatic view of a relational database;

Figure 3 is a highly diagrammatic view of a B-tree index for a relational

database;

5 Figure 4 is a modified dataflow diagram of the method used to calculate DBMS server hardware requirements based on a workload requirement specified explicitly on user entered tpmC;

Figure 5 is a modified dataflow diagram of the method used to calculate DBMS server hardware requirements based on user defined transactions and rates
10 including SQL statement numbers and types including the number of records operated on and selectivity criteria; and

Figure 6 is a flowchart of a process for obtaining any changed hardware utilization limits and making changes in any resulting hardware requirements.

Detailed Description of the Invention

15 Figure 1 illustrates generally a database server system 20 including a server 22 supported by a CRT 24 and a printer 26 for programming, display, maintenance, and general Input/Output uses. Within server 22 is illustrated several CPU sockets 30 and 32, with CPU sockets 30 being populated with CPUs and CPU sockets 32 remaining empty for future expansion and population. Server 22 also includes a memory portion
20 40 which can contain a sufficient quantity of Random Access Memory (RAM) to meet the server's needs. A disk 28 is illustrated for mass storage, which can include disk drives or any other technology capable of holding the contents of the databases or databases to be managed. Several Network Interface Cards (NICs) 42 are illustrated
as part of server 22 and are coupled to a network illustrated by a network link 36
25 which can be any communication link including Local Area Networks, Wide Area

Networks, Ethernet, and the Internet.

Also connected to data link 36 are client computers 38. Software clients can, in fact, reside on the same machine as the server, but in common practice, the client processes usually run on a different machine. In one embodiment, server 22 is a computer running on the Microsoft NT operating system and clients 38 are smaller computers running a Microsoft Windows operating system.

Server 22 is preferably scaleable, having extra socketed capacity for memory, processors, NICs, and disk drives. This allows extra CPUs, memory, NICs, and mass storage such as disk drives to be initially set to meet current needs and later expanded to meet changing needs.

Servers such as server 22 often exist to contain and manage data bases, such as those contained within relational database management systems (RDBMS). RDBMS include tables formed of rows or records and columns. Figure 2 illustrates an RDBMS table 50 formed of several columns 52 and several rows or records 54. Columns 52 typically include both fixed length or width columns and variable length or width columns, where the variable length may be allocated out of a common buffer elsewhere outside of the table itself. In practice, an RDBMS system has numerous tables to be stored and managed.

It is possible for rows 54 to be ordered according to one of the columns. In practice however, the records are typically not ordered, but are linked to indices that are ordered. In a simple example, one of columns 52, such as column 56, may contain a person's social security number and be linked via a software link 58 to an ordered index 60 which contains a sorted list of social security numbers along with the record number at which the corresponding record resides. Such a sorted list of only one column of table 50 can shorten a search from order n to order $\log n$. Such a sorted list

still requires sorting upon the addition, deletion, and change of data.

A commonly used index method is illustrated for column 62, which is linked via a software link 64 to a B-tree index 66. B-tree 66 can contain a multi-level tree well known to those skilled in the software arts. B-tree 66 can be a multi-way tree such as an AVL tree or a 2-3 tree. B-tree indices have the advantage of being quick and easy to modify, without requiring massive amounts of sorting such as in a linear sorted index such as index 60. In particular, a B-tree can be maintained in a balanced condition with the addition of data to avoid skewing the shape of the tree. Maintaining the balance of the tree allows a log n search time to be maintained as well.

In practice, an RDBMS may use only a B-tree for the indexing scheme, due to its utility and flexibility. An RDBMS may maintain a B-tree on any column for which ordered searching may be later requested. As the number of columns to be indexed approaches the number of columns in the table, the data storage requirements for the indices themselves approach and pass the data storage requirements of the table itself. Thus, the data storage requirements of the indices are an important factor to be considered when determining the mass storage requirements for a table and RDBMS.

Figure 3 illustrates a B-Tree 80 including a root node 82 at level 1 having three links 84, 86, and 88 to nodes 90, 92, and 94 respectively at level 2. The nodes at level 2 are illustrated as being doubly linked themselves through links such as links 96 and 98. Links between nodes at the same level, such as links 96 and 98, make maintenance of the B-tree somewhat easier, and browsing can be made somewhat easier as well through use of such links. At level 3, links 100, 102, 104, 106, 108, and 110 are pointed to by the links at level 2. Level 4 is the last level in the tree. B-tree

80 has four levels, or a tree height of four. Level 4 may be said to be the "failure level" of the tree, as level 4 is the level at which a search of the tree will fail if it is to fail. If a value such as a social security number is searched for but there is no such record in the database, level 4 is the level at which the search will fail. At level 4, nodes 112 and 114 are linked together as a doubly linked list by links 113 and 115. In practice, the failure level of a B-Tree is often linked together in this or a similar manner.

In a B-tree, the nodes in the B-tree typically contain only the key or column values the tree is ordered for and points to nodes in the next level related to those keys or column values. For example, in a two-way tree, a node would have one index value, and two pointers, indicating which nodes to go to for values greater than or less than the index value, respectively. B-Trees and databases vary in what they have at the failure level. In some databases, herein termed "physical ordered databases", the failure level has the records themselves linked together. In these databases, once the failure level is reached, the record has been obtained, with no further I/O necessary to obtain the record. In other databases, herein termed "non-physical ordered databases", the nodes at the failure level contain only pointers or record numbers into the main table. In these databases, another I/O is required to obtain the record of interest. In some databases, the failure level contains neither the record of interest nor a pointer to the record of interest. Instead, a unique key is contained, requiring a search on that key to obtain the record of interest. For example, a search of B-Tree ordered on last name plus first name may return only a social security number upon successful completion. Another B-tree or highly optimized index based on social security number can then be rapidly searched for the record of interest. In this scheme, at least one more I/O is required after the failure level has been reached.

number I/Os required to reach a record is of interest to because it determines in part the speed of the database. Both disk I/O and network I/O require latent time to process.

5 In sizing a database, the RDBMS typically has a page size, or an aggregate unit of mass storage typically numbering thousands of bytes. The page size on some computers may be determined in part by operating system and disk system factors. The page size may also be determined by a desire to keep the width of internal variables within the database to manageable limits. The page size is fixed in some RDBMSs and selectable in other RDBMSs.

10 The amount of mass storage required for a single table is a function of several variables, such as the number of rows or records and the number of columns. The database storage required is not a simple calculation of the row size and column sizes for several reasons. First, the column sizes or widths may be variable. Second, the page size enters into the calculation in a non-continuous manner, as some database
15 allocation such as record allocation must lie within a single page rather than cross page boundaries, with some space wasted as a result. Third, some space in a page is set aside for future expansion or reserved for use a buffer space, as when re-ordering data. Fourth, not all space within a page is available for end user data storage, with
20 some being used by the RDBMS itself or for other overhead. In particular, in some RDBMSs, a fraction of each page is specified as not available for initial storage. In some RDBMSs, a number of rows are set aside as non-usable. In some RDBMSs, the size fraction of each record is set aside as non-usable. As previously mentioned, the size of the indices may be a large portion of table storage even though the data itself may not be stored within the indices. All of the aforementioned factors makes sizing the
25 required databases a complicated matter, as is dealt with below.

09545456.022000

Figure 4 illustrates a method 200 for determining the required size of a DBMS server using a transaction throughput benchmark as one of the inputs. In the embodiment illustrated, the TPC-C benchmark is used as the transaction processing environment. The tpmC input refers to the number of New Order Transactions Per Minute for the TPC-C benchmark. In TPC-C, throughput is defined as how many new order transactions per minute a system generates while the system is executing four other transactions types, payment, order-status, delivery, and stock-level; in this benchmark the New Order transaction comprise approximately 40% of the workload.

Method 200 includes inputs 202, outputs 204, and the effective transactions per second (TPS) 206. TPS 206 is the total number of transactions executed per second and is calculated using the equation:

$$\text{TPS} = \text{tpmC} / 36$$

Inputs 202 include the server type 208, the maximum desired processor utilization in percent 210, and the tpmC requirement 212. Optionally, the method may include the tpmC handling ability of a baseline system, indicated at 214, for purposes of comparison. Server type 208 may also be referred to a server family as the server type may include a configurable system, an expandable system, or a family of different hardware systems having varying hardware characteristics. In one embodiment, server type 208 may include varying numbers of processors (CPUs), memory, and mass storage.

The TPC-C benchmark is an Online Transaction Processing (OLTP) benchmark. Key characteristics are a high frequency of transactions executed, a large proportion of the transactions are inserts and updates, and transaction response times should be one second or less. Thus, these characteristics require indexed database access to minimize response time. Further, the database is expected to grow in

proportion to the transaction rate.

A series of measurements were conducted with different numbers of processor configurations and processor speeds. Database sizes were proportionately increased with expected increase in transaction throughput. All benchmark tests were run to 100% processor utilization for the given configuration. The achieved throughput (tpmC), memory used, number of processors, processor speed, mass storage used, number of user simulated were recorded for each benchmark test.

Using the results of these tests a matrix of the above performance and configuration values for each configuration is built. Where configurations are missing interpolation techniques are used to supply the missing values.

For example the following set of published measurement results were used:

System	Mhz	nCPU	tpmC	Mass Storage	No. Users	Memory
Unisys / Aquanta ES2025 Server	550	2	10266	705	8160	1024
Unisys / Aquanta ES2043 Server	500	4	23190	1593	18600	4096
Unisys / Aquanta ES2043R Server	500	4	23190	1593	18600	4096
Unisys / Aquanta ES2045 Server	500	4	23852	1636	19050	4096
Unisys / Aquanta ES5045	500	4	24329	1815	19480	4096
Unisys / Aquanta ES2085R Server	550	8	37757	3079	30480	4096
Unisys / Aquanta ES5085R Server	550	8	40670	3562	32550	4096

This data was then averaged for each configuration and scaled to the lowest Mhz rating. This yielded the following table:

Mhz	nCPU	tpmC	Mass Storage	No. Users	Memory
500	2	9333	641	7418	1024
500	4	23640	1659	18933	4096
500	8	35649	3019	28650	4096

Using curve fitting, this data was then expanded to include estimates at 100 percent processor busy from one to the maximum possible processors. For 500 Mhz processors, the results were as follows:

nCPU	tpmC	Mass Storage	No. Users	Memory
------	------	--------------	-----------	--------

1	6,189	352	4,902	512
2	9,333	641	7,418	1,024
3	16,486	1,150	13,175	1,024
4	23,640	1,659	18,933	4,096
5	26,642	1,999	21,362	4,096
6	29,644	2,339	23,791	4,096
7	32,647	2,679	26,221	4,096
8	35,649	3,019	28,650	4,096

Using scaling techniques, similar results can be calculated for the 550 Mhz processors.

In the method 200 the smallest configuration satisfying both the tpmC 212 requirement and the MaxProcessorUtilization% 210 is selected. The configuration selected then yields the outputs #CPURequired 216 and MemoryRequirements,MB 220. The EffectiveCPUUtilization 218 is calculated as

EffectiveCPUUtilization 218 = (tpmCRequirement 212) / (tpmC at 100% for this configuration). The outputs MassStorageRequirement 220 and #UsersSupported 224 are obtained by interpolating between values for the configurations in the same proportion as tpmC required. For example, if 2 CPUs are required, then the #UsersSupported 224 is calculated as

#UsersSupported = (#UsersSupported with 1 CPU) + (EffectiveCPUUtilization 218) * ((#UsersSupported with 2 CPUs) - (#UsersSupported with 1 CPU))

A similar calculation is made for MassStorageRequirement 222.

~~INS. RE~~ The mass storage required 222 can also be determined by method 200. The mass storage required can be satisfied by adding the appropriate size and number of mass storage devices such as disk drives. The required mass storage can be calculated using the methods described in related applications U.S. Patent Application Serial No. _____, filed _____, entitled ALGORITHM TO CALCULATE MASS STORAGE REQUIREMENTS FOR NT SIZER, and U.S. Patent Application Serial

5 No. _____, filed _____, entitled DATABASE SIZER FOR NT SIZER
SYSTEM, herein incorporated by reference. The tpmC ratio 226 provides a means
for comparing the system specified by the inputs to method 200 and a baseline system
tpmC 214 provided by the user. In one embodiment, a top window in the application
program can be used to selectively browse through the tpmC database. The browsing
can be filtered based on the operating system and database management system. In
one embodiment, the tpmC database is filtered according to operating system and
DBMS to provide a shorter list, which is selectable via a drop down list to a single
tpmC database record, which can be termed the "baseline" system. The tpmC value
10 from the record can be used as the value for the baseline system. In some
embodiments, a known system from one vendor is selected as the baseline system on
one part of the screen. The requirements can be used to select for another system,
termed the 'target' system, from a second vendor in another part of the screen. The
tpmC of the target system can be compared to the tpmC of the baseline system using
15 ratio 226.

Referring now to Figure 5, in another embodiment of the invention, more
detailed inputs are provided to the program to allow direct calculation of the estimated
system configuration requirements. In Figure 5, a method 250 is used to calculate a
workload contribution 258 for each user defined translation. WorkloadContribution
20 258 can be quantified in units of total CPU utilization and NIC utilization. A user
defined transaction can include a transaction name, indicated at TxnName 260, an
expected execution rate of the transaction, indicated at TxnRate 262, a LANSpeed,
indicated at LANSpeed 298, and a transaction composition, indicated at
TxnComposition 264. TxnName 260 can be a user-defined name given to each of the
25 user-defined transactions. TxnRate 262 is the expected execution rate of the user-

defined transaction, typically in transactions per second. TxnComposition 264 includes the user specified SQL composition of the transaction, discussed further below.

5 TxnComposition 264 can be determined by a method 252, and includes the SQL statement compositions 254, including the numbers of those SQL statements, and some parameters included for each SQL statement. SQL statements contributions include insert contributions 268, delete contributions 276, update contributions 278, and select contributions 284, each of which can include the seconds of total CPU time required to execute the SQL and the NIC usage in the total number of bytes passed
10 both ways on the LAN. Method 252 can sum together the workload contributions of each of the SQL statements into total TxnComposition 264. An additional input to method 252 is server type 302 which provides the relative speed with which the SQL will be executed. One method, indicated at 266, calculates the workload contribution of each SQL insert statement, which can be user defined to include a name, indicated
15 at SQLName 270, and the number of identical SQL insert statements, indicated at NoSQLs 272. In general, the larger the number of SQLs, the larger the workload contribution will be.

SQL insert method 266 is typically generates outputs 268 in units of seconds of CPU usage and number of bytes transferred on the LAN. For one CPU usage, a
20 constant number of seconds are required to insert a single record. For LAN usage, the number of bytes transferred is typically a function of the size of the record to be inserted and the average column per row size. For further discussion, see the experimental results section.

SQL delete method 274 calculates the workload contribution of each SQL
25 delete statement, which can be user defined to include a name, indicated at SQLName

270, and the number of identical SQL delete statements, indicated at NoSQLs 272.

SQL delete method 274 is given below and typically generates output 276 in units of seconds of CPU usage and number of bytes transferred on the LAN. For CPU usage, the CPU time required to delete a record is a constant. For LAN usage, the number of bytes transferred is also typically a constant. For further discussion, see the experimental results section.

SQL update method 278 calculates the workload contribution of each SQL update statement, which can be user defined to include a name, indicated at SQLName 270, and the number of records updated for this SQL statement, indicated at NoRecords 288.

SQL update method 278 typically generates output 280 in units of seconds of CPU usage and number of bytes transferred on the LAN. The CPU usage is typically a function of the number of records updated. For LAN usage, the number of bytes transferred is typically a function of the number of records updated. For further discussion, see the experimental results section.

SQL select method 282 calculates the workload contribution of each SQL select statement, which can be user defined to include a name, indicated at SQLName 270, with three other input parameters as the number of tables joined in the select, indicated at NoTables 290; the selectivity, indicated at Selectivity 286; and the number of columns returned, indicated at ColumnsReturned 292.

Select method 282 includes a selectivity criteria input 286, which can vary depending on the type of SQL select statement specified. In a preferred embodiment, the assumption is made that the select statements are indexed which is consistent with OLTP applications. In this same embodiment, for selects involving the joining of two or more tables, the selectivity criteria was geared toward most applicable cases for

OLTP processing. Joins of more than three tables were not included since these cases do not comprise the majority of the transaction activity. Instead, the embodiment concentrated on providing most often occurring cases. To this end, the selectivity was defined for the cited embodiment as follows. If the select is a single table select, selectivity input 286 is the number of rows returned from the table. If the select is a two-table join, then selectivity criteria 286 is the number of rows selected from the outer, or left table, the inner table selectivity is assumed to be an average of four for each row selected from the outer table. If the select is a three-table join, then selectivity 286 is the number of values in the WHERE clauses that pertain to the outer table.

SQL select method 282 typically generates output 284 in units of seconds of CPU usage and number of bytes transferred on the LAN. The CPU usage is typically a complex function of the selectivity, including the table dimensions. For further discussion, see the experimental results section.

For each of the named SQL statements, a series of workload contributions 284, 280, 276, and 268 are generated and summed by method 252 and output as the total TxnComposition 264 which can have units of seconds of CPU usage. TxnComposition 264 can be multiplied by the transaction rate TxnRate 262 to arrive at the transaction workload contribution, TxnWorkloadContribution 258. Each transaction can be named and specified in terms of the SQL statements, and the number of records and/or the selectivity criteria supplied. The total of all TxnWorkLoadContributions 258 can be summed together to arrive at a first estimate of total CPU utilization and total NIC utilization.

The value of the system configuration requirement can be obtained from the method SystemRequirementCalculations 294. The method 294 uses as input the two

requirement can be obtained as a tpmC requirement as discussed with respect to Figure 4. The workload can also be supplied as a more detailed input, as discussed with respect to Figure 5. The workload requirement can also be left at the last value, with no changes. This can be useful when running the program through multiple scenarios. In step 306, the hardware requirements can be calculated while remaining within the hardware utilization limits. In one embodiment, the hardware requirements are determined interactively, with an initial determination of hardware made, followed by a comparison against upper and lower utilization limits. In this embodiment, the percent utilization of some hardware resources such as processors can then be estimated and compared against upper and lower processor utilization limits. If the processor is over utilized then the number of processors can be increased and the computations repeated. If the processor is underutilized, then the number of processors can be decreased and the computations repeated. In another example, if the NICs are over utilized then the number of NICs can be increased and the computations repeated. If the NICs are underutilized, then the number of NICs can be decreased and the computations repeated. In some systems, the speed of the LAN can also be adjusted to provide for increased or decreased NIC utilization.

In step 308, the hardware requirements can be output, typically displayed on a CRT and optionally printed. Hardware requirements can include those discussed with respect to Figure 4, as well as other hardware requirements. In a preferred embodiment, the number of processors used to populate a given server family are output along with a number of NICs to be provided with the system.

In step 310, the user can be prompted for any changes to the utilization limits. In some embodiments, the user can be prompted for changes to workload requirements as well. In either case, if no further changes are desired, the user can

exit as indicated at 312. If changes are desired, changes to utilization limits and/or workload requirements can be made at 314 and new hardware requirements determined. In one example of use, an estimated processor utilization could be viewed on the computer screen and seen to be just above an upper utilization limit.

5 The upper limit could be raised slightly, and the effect on the resulting system configuration viewed. In another example of use, the LAN speed could be adjusted upward, and the effect on the percent utilization and number of NICs viewed. In this way, the number of discrete hardware components required such as NICs and processors may be varied in response to varying workload requirements and varying

10 percent utilization limits.

EXPERIMENTAL RESULTS AND DETAILED METHODS

Measurement Definition and Methodology

A series of 14 SQL statements were defined for measurement on NT Servers running the database products SQL Server 6.5 and Oracle 8.04. These SQL

15 statements, which were run against the tables defined for the 1 GB version of the TPC-D database, were thought to encompass several possible scenarios and business cases for the OLTP environment.

The series of SQL statements, consist of the following types:

- Inserts
- 20 • Deletes
- Updates
- Single table select: three selectivities
- Two table join: three selectivities
- Three table join: three selectivities

The intent behind this set of measurements was that any OLTP transaction can be thought of as some combination of the above SQLs or extensions of them. Thus, the sizer will solicit the user for the series of generic SQLs comprising each transaction. Based on curve fit estimates obtained from the measurement results, one
5 can then estimate the CPU usage for each SQL and subsequently each transaction. By applying workload metrics to each transaction, one can then calculate the CPU utilization and subsequently, the CPU requirement.

The measurement of each SQL was run individually from a client, with no other load on the system. Each test consisted of one or more executions of the same
10 SQL. A SQL may have been run several times in order to obtain a sufficient sample size in the collection of resource usage data. Further, each execution of a given SQL was written to cause the RDBMS to re-parse the SQL before executing it.

Scripts were developed to execute the series of SQLs and to cause certain RDBMS data to be collected and recorded. In order to implement re-parsing for each
15 SQL execution, the script was written accordingly for each RDBMS. Specifically, SQL Server allowed looping and also parsed for each instance within the loop whereas Oracle appeared to parse only once within a loop; consequently, the Oracle SQLs could not be looped, and the scripts were written accordingly.

During each test, performance data was collected from the NT Performance
20 Monitor on the server as well as from the RDBMS data collection construct (SQL Server's ISQL DBCC commands and Oracle's utlbstats/utlestats). This was after the NT Server, the NT Performance Monitor, and the RDBMS were conditioned prior to each test.

Measurement Results

Although several metrics were collected from each measurement, the metrics of most relevance to this version of the sizer were those that measured the CPU resource usage. This consisted of the following:

- 5 • Elapsed time
- Total per cent processor busy
- Per cent processor busy due to the RDBMS

Other statistics, such as logical reads and writes, and physical reads and writes, were also collected. However, this data was not used in the current version of the sizer; this was because all of the table accesses were indexed, and the number of logical IOs was always the number of B Tree levels plus data, consistent with an OLTP environment.

The processor data used in the sizer is the total processor usage, consisting of the RDBMS portion plus the kernel where the kernel activity is assumed to be attributed to the IO activity. In the current version of the sizer the CPU usage due to IO activity is not calculated separately.

Applicant believes that the total CPU time should be approximately the same for a given RDBMS, independent of the number of processors configured. The rationale is that the same amount of work is being done but that the work may be parsed out to a multiple number of processors; the exception is that there is some overhead to parse the work out so that possibly the total CPU time might increase slightly as a function of the number of processors configured. Results from most of the Oracle measurements support this conjecture. It is expected that the same is true for SQL Server

25 Analysis and Curve Fitting

Insert

For inserts, a single record was inserted into a table 10 times followed by a commit. This sequence was repeated 582 times. The CPU time (on a 200 MHz processor system) attributed to inserts, and used in the sizer, was then calculated as

5 $\text{Seconds per record inserted} = (\text{Total CPU Time}) / 582 / 10$

Values used in the sizer are:

SQL Server	.004761
Oracle	.009

Delete

10 The deletes were measured during the process of reverting the database to its original state following the inserts. For the table used, an average of 3.88 records were deleted per delete SQL. In this case, 10 deletes were followed by a commit, and this process was repeated 150 times to return the table to its original condition. The CPU time attributed to deletes, and used in the sizer, was calculated as

15 $\text{Seconds per delete SQL with 3.88 records deleted per SQL} = (\text{Total CPU time}) / 150 / 10$

Values used in the sizer are:

SQL Server	0.015416
Oracle	.015

Update

Three sets of measurements were conducted for update SQLs:

- 20 • A single SQL updated one record; this was followed by a commit
- A single SQL updated five records; this was followed by a commit

Each set of measurements was repeated 100 times.

SQL Server: For the sizer, the CPU time attributed to SQL Server updates was calculated based on the first two points, specifically:

CPU Time = $0.0012491 * (\text{No. records updated per SQL}) + .0004697 \text{ seconds}$

Oracle: The results show fair amount of reduction in overall CPU time for updates as a function of the number of CPUs. The curve fitting for CPU time attributed to Oracle updates was calculated based on both the number of CPUs and the

5 number of records updated per SQL, specifically:

CPU Time = $-0.00248697 * (\text{No. CPUs}) + 0.00033591 * (\text{No. records updated per SQL}) + 0.00995405 \text{ seconds}$

Currently the sizer pessimistically uses calculations based on 1 CPU.

Single Table Select - Indexed

10 Three single table SELECT SQLs were individually run. The WHERE clauses were set to produce selectivities of 2,588, 12,477, 24,896 records, respectively. These SQLs were defined so that the access method of the RDBMS was indexed, consistent with an OLTP environment. For SQL Server this was via a cluster key, thus, eliminating the additional IO to another leaf page if instead the
15 index were non-clustered. For Oracle, this was via a standard index. The most suitable curve fit for both RDBMS seemed to be to divide the CPU time by the number of records selected, and then averaging that value. A value of .1 msec was given each RDBMS for parsing the SQL

The resulting formulas used, and applied in the sizer, were the following:

20 SQL Server: $(\text{CPU Time}) = .00006.73E * (\text{Records Selected}) + .0001$

Oracle: $(\text{CPU Time}) = 0.00218 * (\text{Records Selected}) + .0001$

Note that the usual methods of curve fitting would apply to the selectivities noted, i.e., a very large number of records selected. Extrapolating to a few records selectivity (consistent with OLTP) produced either extremely large or extremely small

CPU time overheads. As a result, the method of averaging time per record selected was chosen as a means of estimation.

Two Table Join - Indexed

For the two table join, an example of the business basis for the SQL SELECT was to find the various sources (e.g., prices) of a given item or list of items. In this case, there were on the average four such sources for each item. Three SQL SELECTs were constructed to return the sources for one, five, and 10 items, respectively. Results of the measurements for each RDBMS are shown in the charts below.

SQL Server: Curve fitting based on the template $y = a * x^b$ produced the following results:

$$(\text{CPU Time}) = .0014665 * (\text{Items Selected}) ^ .7552$$

Oracle: Curve fitting based on the template $y = a * x + b$ produced the following results:

$$(\text{CPU Time}) = 0.00858 * (\text{Items Selected}) + 0.003236$$

Three Table Join - Indexed

For the three table join, an example of the business basis for the SQL SELECT was to determine the status of orders placed on certain dates from a selected segment of the customer population. For this particular database, the customer segment chosen places about 20% of all of the orders. Each order consists of approximately 4 items on the average. For this set of measurements three SQLs were defined to return status on one, five, and 10 order dates.

Curve fitting for both the SQL Server and Oracle RDBMS' was based on the template $y = a * x + b$. For the Oracle case, the measurement results were first

averaged over the number of CPUs for each point; then linear curve fitting was applied to the averages.

The resulting formulas used, and applied in the sizer, are the following:

SQL Server: $(\text{CPU Time}) = 0.001306 * (\text{Selectivity on Outer Table}) -$

5 0.00039

Oracle: $(\text{CPU Time}) = 0.002521 * (\text{Selectivity on Outer Table}) + 0.001$

LAN USAGE

This section provides the formulas used to calculate the LAN bytes passed per SQL statement.

10 Using the Ethernet and TCP/IP protocol, a certain amount of bytes are passed with each message. Some of these values are as follows:

FrameHeader = 72 bytes

When a SQL arrives or when it is either simply acknowledged or data is passed back, the following assumptions are made as to the amount of data transferred

15 across the LAN:

SQLInsertIncomingData = 200 bytes

SQLUpdateIncomingData = 200 bytes

SQLDeleteIncomingData = 200 bytes

SQLSelectIncomingData = 200 bytes

20 SQLAcknowledge = 100 bytes

The amount of data CommBytes transferred both ways across the LAN is estimated as follows per SQL:

Inserts:

25 $\text{CommBytes} = 2 * \text{FrameHeader} + \text{SQLInsertIncomingData} +$
 $\text{AverageColumnsPerRow} + \text{SQLAcknowledge}$
Deletes:

CommBytes = 2 * FrameHeader + SQLDeleteIncomingData +
SQLAcknowledge

Updates:

CommBytes = 2 * FrameHeader + SQLUpdateIncomingData * NoRecords +
5 SQLAcknowledge
Selects:

1 Table: Selected Bytes = Selectivity * BytesPerColumn * NoColumns

2 Tables: Selected Bytes = 4 * Selectivity * BytesPerColumn * NoColumns

3 Tables: Selected Bytes = 5 * Selectivity * BytesPerColumn * NoColumns

10 CommBytes = FrameHeader + SQLSelectIncomingData

Do While SelectedBytes > 1500 ' too big for a frame, break up into 1500

bytes

CommBytes = CommBytes + TCPIPHdr + 1500

SelectedBytes = SelectedBytes - 1500

15 Loop

If SelectedBytes > 0 Then

CommBytes = CommBytes + TCPIPHdr + SelectedBytes

End If

20

HEADROOM CALCULATIONS

This SECTION describes a method for calculating the system configuration requirement and the headroom.

Processor Requirement

25 Define

TxnRate_I = transaction rate of transaction I

CPUTime_I = total CPUTime used to execute one instance of transaction I

$TxnRate_T = \text{sum of } TxnRate_I \text{ over all transactions}$

For each transaction, an input to method 294 is the transaction's effective contribution to CPU utilization, expressed as

$$(CPUTime_I) * (TxnRate_I)$$

5 Define

$Util_1CPU = \text{sum of } (CPUTime_I) * (TxnRate_I) \text{ over all transactions}$

So $Util_1CPU$ is defined as the effective utilization of a single CPU. Note that this value can exceed 100%, implying that more than one CPU is required to service the transaction rate.

10 One method to determine the CPU requirement NoCPUs is as follows:

$$NoCPUs = \text{CEILING}(Util_1CPU/100)$$

While $(Util_1CPU/100)/NoCPUs > \text{MaxProcessorUtilization}$

$$NoCPUs = NoCPUs + 1$$

Wend

15 We note, however, that adding multiple processors does not increase performance with 100% efficiency. Thus, for example, by doubling the number of processors, we achieve anywhere from 150% to 180% as much performance, rather than 200%. Much of this is due to a combination of the instruction speed slows due to locality of reference and the number of instructions per transaction increase due to
20 increased housekeeping overhead.

The service time per transaction is the quotient (instructions per transaction) / (instructions per second). It has been found in studies that the service time per transaction is a linear function of the number of processors, i.e., service time per transaction, n processors = (service time, 1 processor) * (a * n + b) where a and b are

constants determined from analysis of measurement results. Thus, the transaction rate is proportional to the reciprocal, that is,
 transaction rate $\sim 1 / (a * n + b)$

For this version of the sizer, the values a, b are:

5 $a = 0.028982$

$b = 0.973703$

The above algorithm is then modified as follows:

Define

$NoCPUs = CEILING(Util_1CPU/100)$

10 $CPUSvc_1x = Util_1CPU / Txn_Rate_T$

$bResult = False$

While Not $bResult$

$temp = NoCPUs / (a * NoCPUs + b) * MaxCPUUtilization$

If $temp < Util_1CPU$ Then

15 $NoCPUs = NoCPUs + 1$

Else

$bResult = True$

$EffectiveCPUUtilization = Ceiling((a * NoCPUs + b) * Util_1CPU / NoCPU * 100)$

End If

20 Wend

Using this method, if the selected processor type does not meet the transaction method, then the next faster processor is selected, and the above process is repeated.

The headroom algorithm ensures the following:

LAN speeds."

End If

This algorithm also demonstrates the headroom algorithm.

Memory Calculations

5 Because the database servers are of the enterprise quality, the amount of memory required is generally very high. As a rule of thumb the memory requirement has been specified as an additional 512 MB per processor up to the limit imposed by the machine and the operating system.

10 Numerous advantages of the invention covered by this document have been set forth in the foregoing description. It will be understood, however, that this disclosure is, in many respects, only illustrative. Changes may be made in details, particularly in matters of shape, size, and arrangement of parts without exceeding the scope of the invention. The invention's scope is, of course, defined in the language in which the appended claims are expressed.